



Variable selection in model-based discriminant analysis

Cathy Maugis, Gilles Celeux, Marie-Laure Martin-Magniette

► To cite this version:

Cathy Maugis, Gilles Celeux, Marie-Laure Martin-Magniette. Variable selection in model-based discriminant analysis. [Research Report] RR-7290, INRIA. 2010. inria-00483229

HAL Id: inria-00483229

<https://hal.inria.fr/inria-00483229>

Submitted on 12 May 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

Variable selection in model-based discriminant analysis

Cathy Maugis — Gilles Celeux — Marie-Laure Martin-Magniette

N° 7290

May 2010

Thème COG

A large, light gray stylized 'R' logo that serves as a background for the text.

*Rapport
de recherche*

Variable selection in model-based discriminant analysis

Cathy Maugis ^{*}, Gilles Celeux [†], Marie-Laure Martin-Magniette^{‡§}

Thème COG — Systèmes cognitifs
Équipe-Projet SELECT

Rapport de recherche n° 7290 — May 2010 — 31 pages

Abstract: A general methodology for selecting predictors for Gaussian generative classification models is presented. The problem is regarded as a model selection problem. Three different roles for each possible predictor are considered: a variable can be a relevant classification predictor or not, and the irrelevant classification variables can be linearly dependent on a part of the relevant predictors or independent variables. This variable selection model was inspired by the model-based clustering model of Maugis et al. (2009b). A BIC-like model selection criterion is proposed. It is optimized through two embedded forward stepwise variable selection algorithms for classification and linear regression. The model identifiability and the consistency of the variable selection criterion are proved. Numerical experiments on simulated and real data sets illustrate the interest of this variable selection methodology. In particular, it is shown that this well ground variable selection model can be of great interest to improve the classification performance of the quadratic discriminant analysis in a high dimension context.

Key-words: Discriminant, redundant or independent variables, Variable selection, Gaussian classification models, Linear regression, BIC

^{*} Institut de Mathématiques de Toulouse, INSA de Toulouse, Université de Toulouse

[†] INRIA Saclay - Île-de-France, Projet SELECT, Université Paris-Sud 11

[‡] UMR AgroParisTech/INRA MIA 518, Paris

[§] URGV UMR INRA 1165, UEVE, ERL CNRS 8196, Evry

Sélection de variables pour l'analyse discriminante gaussienne

Résumé : Nous proposons une méthodologie générale pour la sélection de variables en analyse discriminante par des modèles génératifs gaussiens. Le problème est vu sous un angle de choix de modèles. Les variables en compétition peuvent avoir trois rôles : ce sont soit des prédicteurs utiles pour la classification supervisée, soit des variables redondantes, liés aux prédicteurs par une régression linéaire, soit des variables indépendantes. Ce modèle s'inspire directement du modèle de Maugis et al. (2009b) pour la sélection de variables en classification non supervisée par des modèles de mélanges de lois gaussiennes. Un critère de type BIC est proposé pour choisir le rôle des variables. Ce critère est optimisé par deux algorithmes emboîtés de sélection ascendante avec remise en cause pour la classification et la régression. Nous établissons l'identifiabilité de notre modèle et nous prouvons l'optimalité asymptotique de notre critère. Nous illustrons les bonnes performances de notre approche par des expérimentations sur des données simulées et réelles. Nous montrons en particulier que notre méthodologie de sélection de variables peut être profitable pour l'analyse discriminante quadratique en grand dimension.

Mots-clés : Variables discriminantes, redondantes ou indépendantes, Sélection de variables, Classification supervisée gaussienne, Régression linéaire, BIC

1 Introduction

The task of supervised classification is to build a classifier which enables us to assign an object described by predictors to one of known classes. Such classifiers are built from a training set of objects for which the predictor measurements and the class labels are known. A lot of different methods are available, see for instance the recent books on statistical learning by Hastie et al. (2009) or Bishop (2006). Those methods differ in the way they approach the problem. Generative models, as Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA), estimate the class-conditional densities. The predictive models (for instance logistic regression, classification trees and the k -nearest-neighbor classifier) directly estimate the posterior class probabilities. Non probabilistic methods, as Neural networks and Kernel methods such as Support Vectors, aim at finding the decision function which characterizes the classifier.

Generative models are less parsimonious than predictive and non probabilistic methods. Those last methods are generally preferred to generative models when the number of predictors is large in regard to the number of objects in the training set. However, generative models have some advantages since they allow us to determine the marginal density of the data. As noted in Hastie et al. (2009), LDA and QDA are widely used and perform well on an amazingly large and diverse set of classification problems. Moreover LDA is regarded as a reference method by many practitioners, and an advantage of LDA over QDA is that it is a more parsimonious method.

Much efforts have been paid in variable selection for classification, see the reviews of Guyon and Elisseeff (2003) and Mary-Huard et al. (2007). In this paper, we concentrate our attention on variable selection for Gaussian generative models. There exists quite efficient methods to select predictors in the LDA context. Efficient stepwise variable selection procedures are available in most statistical softwares (see McLachlan, 1992, Section 12.3.3). On the contrary, there is less available material for QDA (Young and Odell, 1986), and as far as we know, no variable selection procedure for QDA is available in standard statistical softwares. However in the last few years, there is a renewal of interest in this topic. Zhang and Wang (2008) proposed a variable selection procedure for QDA based on a BIC criterion and Murphy et al. (2010) have adapted the variable selection procedure of Raftery and Dean (2006) to the supervised classification context.

The purpose of this paper is to extend the general variable selection modelling proposed in Maugis et al. (2009b), conceived for model-based clustering to the Gaussian classification models. This modelling is the result of successive improvements of variable selection modelling in model-based clustering (Raftery and Dean, 2006; Maugis et al., 2009a,b). Acting in such a way, we dramatically strengthen the appeal of non linear Gaussian classifiers, proposed by Bensmail and Celeux (1996) which are up to now limited by the large number of parameters to be estimated. The models and variable selection algorithms proposed not only lead to interpret the roles of variables in a clear way, but they also lead to much increase the discriminative efficiency of methods such as QDA.

The paper is organized as follows. In Section 2, Gaussian models of classification are recalled. Our variable selection approach is presented in Section 3. It makes use of a model which states a clear distinction between useful, redundant and noisy variables for the classification task in the Gaussian framework. It

leads to a BIC-like criterion to be optimized (Section 4). It is proved in Section 5 that our approach leads to identifiable classification models and that our variable selection criterion is consistent under mild assumptions. In Section 6, a variable selection algorithm using two forward stepwise algorithms is described to determine the roles of the predictors. Applications on simulated and real data sets are presented in Section 7. A short discussion section ends the paper and the proofs of the theorems of Section 5 are postponed to Appendices B and C.

2 Gaussian classification models

Training data for discriminant analysis are composed by n vectors

$$(\underline{\mathbf{x}}, \underline{z}) = \{(\mathbf{x}_1, z_1), \dots, (\mathbf{x}_n, z_n); \mathbf{x}_i \in \mathbb{R}^Q, z_i \in \{1, \dots, K\}\},$$

where \mathbf{x}_i is the Q -dimensional predictor and z_i is the class label of the i th subject. We assume that the prior probability of the class G_k is $P(z = k) = p_k$ with $p_k > 0$ for any k , $1 \leq k \leq K$ and $\sum_{k=1}^K p_k = 1$. The class conditional density of class G_k is modelled with a Q -dimensional Gaussian density: $\mathbf{x}_i | z_i = k \sim \mathcal{N}_Q(\mu_k, \Sigma_k)$ where $\mu_k \in \mathbb{R}^Q$ is the mean vector and Σ_k is the $Q \times Q$ variance matrix. The aim of discriminant analysis is to design a classifier from the training sample, allowing us to estimate the label of any new observation $\mathbf{x} \in \mathbb{R}^Q$.

Gaussian generative models differ essentially in their assumptions about the variance matrices. The most commonly applied method, called linear discriminant analysis (LDA), assumes that the variance matrices of the different class are equal. When the variance matrices are totally free, the method is called quadratic discriminant analysis (QDA). Bensmail and Celeux (1996) generalize the LDA and QDA methods in the Eigenvalue Decomposition Discriminant Analysis (EDDA). As in Banfield and Raftery (1993) and Celeux and Govaert (1995), EDDA is based on the eigenvalue decomposition of the variance matrices

$$\forall k \in \{1, \dots, K\}, \Sigma_k = L_k D_k A_k D_k'$$

where $L_k = |\Sigma_k|^{\frac{1}{Q}}$, D_k is the Σ_k 's eigenvector matrix and A_k is the diagonal matrix of the normalized eigenvalues of Σ_k . Those elements respectively control the volume, the orientation and the shape of the density contour of class G_k . According to constraints required on the three elements of the eigenvalue decomposition, a collection \mathcal{M} of 14 more or less parsimonious and easily interpreted models is available (see Table 4 in Appendix A). Those 14 models are available in the MIXMOD software (Biernacki et al., 2006) and, for most of them, in the MCLUST software (Fraley and Raftery, 2003). The LDA and QDA are implemented in several softwares as well as R in the library MASS.

The model selection in the EDDA context consists of choosing the best form of the variance matrices. The best model is usually selected by minimising the cross-validated classification error rate (Bensmail and Celeux, 1996). Another possible selection criterion is the Bayesian information criterion (Schwarz, 1978) which is an asymptotic approximation of the integrated loglikelihood. The model which maximizes the Bayesian information criterion (BIC) is selected. In this paper a selection by the BIC is considered. Notice that BIC focuses on the

model fit rather than the minimisation of the misclassification rate and is much cheaper to compute.

Once a model of the collection is selected, a new observation \mathbf{x}_0 is assigned to the group for which its a posteriori probability is maximum. It is the *Maximum A Posteriori* rule and it is equivalent to find the class k^* such that

$$k^* = \operatorname{argmax}_{1 \leq k \leq K} p_k \Phi(\mathbf{x} | \mu_k, \Sigma_{k(m)}),$$

where $\Phi(\cdot | \mu_k, \Sigma_{k(m)})$ denotes the Gaussian density with mean vector μ_k and variance matrix $\Sigma_{k(m)}$ fulfilling the form $m \in \mathcal{M}$.

3 The variable selection model collection

Each of the Q available variables brings information (its own ability to separate the classes), and noise (its sampling variance). Thus it is important to select the variables bringing more discriminant information than noise. In practice there are three kinds of variables: The discriminant variables useful for the classification task, the redundant variables linked to the discriminant variables, and the noisy variables which bring no information for the classification task. Thus, variable selection is an important part of discriminant analysis to get a reliable and parsimonious classifier. Considering the classification problem in the model-based discriminant analysis context allows us to recast variable selection into a model selection problem and to adapt the variable selection model for model-based clustering of Maugis et al. (2009b) in the supervised classification context.

In our modelling, the variables have three possible roles: relevant, redundant or independent for the discriminant analysis. The nonempty set of relevant predictors is denoted as S and the independent variable subset is denoted as W . The redundant variables, whose the subset is denoted as U , are explained by a variable subset R of S according to a linear regression while the variables in W are assumed to be independent of all the relevant variables. Note that if U is empty, R is empty too and otherwise R is assumed to be not empty. Thus denoting \mathcal{F} the family of variable index subsets of $\{1, \dots, Q\}$, the variable partition set can be described as follows:

$$\mathcal{V} = \left\{ (S, R, U, W) \in \mathcal{F}^4; \begin{array}{l} S \cup U \cup W = \{1, \dots, Q\} \\ S \cap U = \emptyset, S \cap W = \emptyset, U \cap W = \emptyset \\ S \neq \emptyset, R \subseteq S \\ R = \emptyset \text{ if } U = \emptyset \text{ and } R \neq \emptyset \text{ otherwise} \end{array} \right\}.$$

Throughout this paper, a quadruplet (S, R, U, W) of \mathcal{V} is denoted as $\mathbf{V} = (S, R, U, W)$.

The law of the training sample is modelled by, $\forall(\mathbf{x}, z) \in \mathbb{R}^Q \times \{1, \dots, K\}$,

$$\begin{cases} f(\mathbf{x} | z = k, m, r, l, \mathbf{V}) &= \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)}) \Phi(\mathbf{x}^W | \gamma, \tau_{(l)}) \\ (\mathbb{I}_{z=1}, \dots, \mathbb{I}_{z=K}) &\sim \text{Multinomial}(1; p_1, \dots, p_K) \end{cases}$$

where

- on the discriminant variable subset S , the variance matrices $\Sigma_{1(m)}, \dots, \Sigma_{K(m)}$ fulfill the constraints of the form $m \in \mathcal{M}$ (see Section 2);
- on the redundant variable subset U , the density $\Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)})$ corresponds to the linear regression density of \mathbf{x}^U on \mathbf{x}^R , where the vector a is the intercept vector, β is the regression coefficient matrix and $\Omega_{(r)}$ is the variance matrix; this last matrix is assumed to have a spherical ($[LI]$), diagonal ($[LB]$) or a general ($[LC]$) form, and this form is specified by $r \in \mathcal{T}_{reg} = \{[LI], [LB], [LC]\}$;
- on the independent variable subset W , the marginal density is assumed to be a Gaussian density with mean γ and variance matrix $\tau_{(l)}$ which can be spherical or diagonal and is specified by $l \in \mathcal{T}_{indep} = \{[LI], [LB]\}$.

Finally the model collection is

$$\mathcal{N} = \{(m, r, l, \mathbf{V}); m \in \mathcal{M}, r \in \mathcal{T}_{reg}, l \in \mathcal{T}_{indep}, \mathbf{V} \in \mathcal{V}\} \quad (1)$$

and the likelihood of model (m, r, l, \mathbf{V}) is given by

$$f(\mathbf{z}, \mathbf{z} | m, r, l, \mathbf{V}, \theta) = \prod_{i=1}^n \prod_{k=1}^K [p_k \Phi(\mathbf{x}_i^S | \mu_k, \Sigma_{k(m)}) \Phi(\mathbf{x}_i^U | a + \mathbf{x}_i^R \beta, \Omega_{(r)}) \Phi(\mathbf{x}_i^W | \gamma, \tau_{(l)})]^{\mathbf{1}_{z_i=k}}$$

where the parameter vector $\theta = (\alpha_{(m)}, a, \beta, \Omega_{(r)}, \gamma, \tau_{(l)})$ with

$$\alpha_{(m)} = (p_1, \dots, p_K, \mu_1, \dots, \mu_K, \Sigma_{1(m)}, \dots, \Sigma_{K(m)})$$

belongs to a parameter vector set $\Upsilon_{(m,r,l,\mathbf{V})}$.

4 Model selection criterion

The model collection \mathcal{N} allows us to recast the variable selection problem for Gaussian discriminant analysis into a model selection problem. Ideally, we search the model maximizing the integrated loglikelihood

$$(\tilde{m}, \tilde{r}, \tilde{l}, \tilde{\mathbf{V}}) = \underset{(m,r,l,\mathbf{V}) \in \mathcal{N}}{\operatorname{argmax}} \ln[f(\mathbf{z}, \mathbf{z} | m, r, l, \mathbf{V})]$$

where

$$f(\mathbf{z}, \mathbf{z} | m, r, l, \mathbf{V}) = \int f(\mathbf{z}, \mathbf{z} | m, r, l, \mathbf{V}, \theta) \Pi(\theta | m, r, l, \mathbf{V}) d\theta,$$

Π being the prior distribution of the vector parameter. Since this integrated loglikelihood is difficult to evaluate, it could be approximated by the BIC criterion (Schwarz, 1978). Then the selected model satisfies

$$(\hat{m}, \hat{r}, \hat{l}, \hat{\mathbf{V}}) = \underset{(m,r,l,\mathbf{V}) \in \mathcal{N}}{\operatorname{argmax}} \operatorname{crit}(m, r, l, \mathbf{V}) \quad (2)$$

where the model selection criterion is defined by

$$\operatorname{crit}(m, r, l, \mathbf{V}) = \operatorname{BIC}_{da}(\mathbf{x}^S, \mathbf{z} | m) + \operatorname{BIC}_{reg}(\mathbf{x}^U | r, \mathbf{x}^R) + \operatorname{BIC}_{indep}(\mathbf{x}^W | l), \quad (3)$$

where

- the BIC criterion for the Gaussian discriminant analysis on the relevant variable subset S is given by

$$\text{BIC}_{\text{da}}(\underline{\mathbf{x}}^S, \underline{z}|m) = 2 \sum_{i=1}^n \ln \left[\sum_{k=1}^K \hat{p}_k \Phi(\mathbf{x}_i^S | \hat{\mu}_k, \hat{\Sigma}_{k(m)}) \mathbb{I}_{z_i=k} \right] - \lambda_{(m,S)} \ln(n)$$

where $\hat{\alpha}_{(m)}$ is the maximum likelihood estimator and $\lambda_{(m,S)}$ is the number of free parameters for the model m on the variable subset S .

- the BIC criterion for the linear regression of the variable subset U on R is defined by

$$\text{BIC}_{\text{reg}}(\underline{\mathbf{x}}^U | r, \underline{\mathbf{x}}^R) = 2 \sum_{i=1}^n \ln[\Phi(\mathbf{x}_i^U | \hat{a} + \mathbf{x}_i^R \hat{\beta}, \hat{\Omega}_{(r)})] - \nu_{(r,U,R)} \ln(n) \quad (4)$$

where \hat{a} , $\hat{\beta}$ and $\hat{\Omega}_{(r)}$ are the maximum likelihood estimators and $\nu_{(r,U,R)}$ is the number of free parameters of the linear regression.

- the BIC criterion associated to the Gaussian density on the variable subset W is given by

$$\text{BIC}_{\text{indep}}(\underline{\mathbf{x}}^W | l) = 2 \sum_{i=1}^n \ln[\Phi(\mathbf{x}_i^W | \hat{\gamma}, \hat{\tau}_{(l)})] - \rho_{(l,W)} \ln(n).$$

The parameters $\hat{\gamma}$ and $\hat{\tau}_{(l)}$ denote the maximum likelihood estimators and $\rho_{(l,W)}$ is the number of free parameters of the Gaussian density.

- the maximum likelihood estimator is denoted as $\hat{\theta} = (\hat{\alpha}_{(m)}, \hat{a}, \hat{\beta}, \hat{\Omega}_{(r)}, \hat{\gamma}, \hat{\tau}_{(l)})$ and the overall number of free parameters is $\Xi_{(m,r,l,\mathbf{v})} = \lambda_{(m,S)} + \nu_{(r,U,R)} + \rho_{(l,W)}$.

5 Theoretical properties

The theoretical properties established in Maugis et al. (2009b) in the model-based clustering framework can be adapted to the Gaussian discriminant analysis context. First, necessary and sufficient conditions are given to ensure the identifiability of the model collection. Second, a consistency theorem of the model selection criterion is stated.

5.1 Identifiability

In order to ensure the model identifiability, some natural conditions are required to distinguish the discriminant density part to the regression and the independent Gaussian density parts. For instance, if s is a non empty subset strictly included into the relevant variable subset S and \bar{s} is its complement in S then the identifiability cannot be ensured if the regression density of \bar{s} on s can be regrouped with the regression density of U on R . Despite the fact that Conditions (C1)-(C3) of Theorem 1 look rather technical, they are quite natural and Theorem 1 is saying that our variable selection model is identifiable in all situations of interest.

The following additional notation is introduced to state the model identifiability theorem. Recall that $\Phi(\cdot|\mu_k, \Sigma_k)$ denotes the Gaussian density with mean μ_k and variance matrix Σ_k . The parameters can be decomposed into $\mu_k = (\mu_{ks}, \mu_{k\bar{s}})$ and Σ_k into submatrices $\Sigma_{k,ss}$, $\Sigma_{k,s\bar{s}}$ and $\Sigma_{k,\bar{s}\bar{s}}$, where s is a nonempty subset of S and \bar{s} its complement in S . Moreover, conditional parameters are defined by $\mu_{k,\bar{s}|s} = \mu_{k\bar{s}} - \mu_{ks}\Sigma_{k,ss}^{-1}\Sigma_{k,s\bar{s}}$, $\Sigma_{k,\bar{s}|s} = \Sigma_{k,\bar{s}\bar{s}} - \Sigma_{k,s\bar{s}}\Sigma_{k,ss}^{-1}\Sigma_{k,s\bar{s}}$ and $\Sigma_{k,\bar{s}\bar{s}|s} = \Sigma_{k,\bar{s}\bar{s}} - \Sigma_{k,s\bar{s}}\Sigma_{k,ss}^{-1}\Sigma_{k,s\bar{s}}$. For two subsets s and t , the following restrictions of a $I \times J$ matrix Λ are considered: $\Lambda_{st} = (\Lambda_{ij})_{i \in s, j \in t}$, $\Lambda_{.t} = (\Lambda_{ij})_{1 \leq i \leq I, j \in t}$ and $\Lambda_{.s} = (\Lambda_{ij})_{i \in s, 1 \leq j \leq J}$.

Theorem 1. *Let $\Theta_{(m,r,l,\mathbf{V})}$ be a subset of the parameter set $\Upsilon_{(m,r,l,\mathbf{V})}$ such that elements $\theta = (\alpha, a, \beta, \Omega, \gamma, \tau)$*

(C1) : *contain couples (μ_k, Σ_k) fulfilling $\forall s \subsetneq S, \exists (k, k'), 1 \leq k < k' \leq K$*

$$\mu_{k,\bar{s}|s} \neq \mu_{k',\bar{s}|s} \text{ or } \Sigma_{k,\bar{s}|s} \neq \Sigma_{k',\bar{s}|s} \text{ or } \Sigma_{k,\bar{s}\bar{s}|s} \neq \Sigma_{k',\bar{s}\bar{s}|s},$$

where \bar{s} denotes the complement in S of any nonempty subset s of S .

(C2) : *if $U \neq \emptyset$,*

- * for all variables j of R , there exists a variable u of U such that the restriction β_{uj} of the regression coefficient matrix β associated with j and u is not equal to zero.*
- * for all variables u of U , there exists a variable j of R such that $\beta_{uj} \neq 0$.*

(C3) : *Parameters Ω and τ strictly respect the forms r and l respectively: They are both diagonal matrices with at least two different eigenvalues if $r = [LB]$ and $l = [LB]$ and Ω has at least a non-zero entry outside the main diagonal if $r = [LC]$.*

Let (m, r, l, \mathbf{V}) and $(m^, r^*, l^*, \mathbf{V}^*)$ be two models. If there exist $\theta \in \Theta_{(m,r,l,\mathbf{V})}$ and $\theta^* \in \Theta_{(m^*,r^*,l^*,\mathbf{V}^*)}$ such that*

$$f(\cdot|m, r, l, \mathbf{V}, \theta) = f(\cdot|m^*, r^*, l^*, \mathbf{V}^*, \theta^*)$$

then $(m, r, l, \mathbf{V}) = (m^, r^*, l^*, \mathbf{V}^*)$ and $\theta = \theta^*$.*

The complete proof of Theorem 1 is postponed to Appendix B.

5.2 Consistency

A consistency property of our criterion can be checked. In this section, it is proved that the probability of selecting the true model by maximizing Criterion (3) approaches 1 as $n \rightarrow \infty$. Denoting h the density function of the sample (\mathbf{x}, \mathbf{z}) , the two following vectors are considered

$$\begin{aligned} \theta_{(m,r,l,\mathbf{V})}^* &= \underset{\theta_{(m,r,l,\mathbf{V})} \in \Theta_{(m,r,l,\mathbf{V})}}{\operatorname{argmin}} \quad \text{KL}[h, f(\cdot|m, r, l, \mathbf{V}, \theta)] \\ &= \underset{\theta_{(m,r,l,\mathbf{V})} \in \Theta_{(m,r,l,\mathbf{V})}}{\operatorname{argmax}} \quad \mathbb{E}\{\ln f(X, Z|m, r, l, \mathbf{V}, \theta)\}, \end{aligned}$$

where $\text{KL}[h, f] = \int \ln \left\{ \frac{h(x)}{f(x)} \right\} h(x) dx$ is the Kullback-Leibler divergence between the densities h and f and

$$\hat{\theta}_{(m,r,l,\mathbf{V})} = \underset{\theta_{(m,r,l,\mathbf{V})} \in \Theta_{(m,r,l,\mathbf{V})}}{\text{argmax}} \quad \frac{1}{n} \sum_{i=1}^n \ln \{f(\mathbf{x}_i, z_i | m, r, l, \mathbf{V}, \theta)\}.$$

Recall that $\Theta_{(m,r,l,\mathbf{V})}$'s are the subsets defined in Theorem 1 for ensuring the model identifiability.

The following assumption is considered:

- (H1) The density h is assumed to be one of the densities in competition. By identifiability, there exists a unique model $(m_0, r_0, l_0, \mathbf{V}_0)$ and an associated parameter θ^* such that $h = f(\cdot | m_0, r_0, l_0, \mathbf{V}_0, \theta^*)$.

Moreover, an additional technical assumption is considered:

- (H2) For all models $(m, r, l, \mathbf{V}) \in \mathcal{N}$, the vectors θ^* and $\hat{\theta}$ are supposed to belong to a compact subspace $\Theta'_{(m,r,l,\mathbf{V})}$ in the intersection between $\Theta_{(m,r,l,\mathbf{V})}$ and

$$\left(\begin{array}{c} \mathcal{P}_{K-1}(\rho) \times \mathcal{B}(\eta, \text{card}(S))^K \times \mathcal{D}_{\text{card}(S)}^K \times \mathcal{B}(\eta, \text{card}(U)) \\ \times \mathcal{B}(\eta, \text{card}(R), \text{card}(U)) \times \mathcal{D}_{\text{card}(U)} \times \mathcal{B}(\eta, \text{card}(W)) \times \mathcal{D}_{\text{card}(W)} \end{array} \right)$$

where

- $\mathcal{P}_{K-1}(\rho) = \left\{ (p_1, \dots, p_K) \in [\rho, 1]^K; \sum_{k=1}^K p_k = 1 \right\}$ where $\rho > 0$,
- $\mathcal{B}(\eta, r)$ is the closed ball in \mathbb{R}^r of radius η centered at zero for the l^2 -norm defined by $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^r x_i^2}$, $\forall \mathbf{x} \in \mathbb{R}^r$,
- $\mathcal{B}(\eta, r, q)$ is the closed ball in $\mathcal{M}_{r \times q}(\mathbb{R})$ of radius η centered at zero for the matricial norm $\|\cdot\|$ defined by

$$\forall A \in \mathcal{M}_{r \times q}(\mathbb{R}), \|A\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{x}A\|,$$

- \mathcal{D}_r is the set of the $r \times r$ positive definite matrices with eigenvalues in $[s_m, s_M]$ with $0 < s_m < s_M$.

Theorem 2. Under assumptions (H1) and (H2), the model $(\hat{m}, \hat{r}, \hat{l}, \hat{\mathbf{V}})$ maximizing Criterion (3) is such that

$$P((\hat{m}, \hat{r}, \hat{l}, \hat{\mathbf{V}}) = (m_0, r_0, l_0, \mathbf{V}_0)) \xrightarrow{n \rightarrow \infty} 1.$$

The proof of this theorem is given in Appendix C.

6 The variable selection procedure

Theorem 2 is reassuring about the theoretical behavior of the model selection Criterion (3). Unfortunately, the number of models given by (1) being huge, an exhaustive search for the model maximizing Criterion (3) is impossible. Thus we design a procedure, embedding forward stepwise algorithms, to determine the best variable roles and the best variance matrix forms.

6.1 The models in competition

At a fixed step of the algorithm, the variable set $\{1, \dots, Q\}$ is divided into the subset of selected discriminant variables S , the subset U of redundant variables which are linked to some discriminant variables, the subset W of independent irrelevant variables and j the candidate variable for inclusion into or exclusion from the discriminant variable subset. Under the model (m, r, l) , the integrated likelihood can be decomposed as

$$\begin{aligned} f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{\mathbf{x}}^U, \underline{\mathbf{x}}^W, \underline{z}|m, r, l) &= f(\underline{\mathbf{x}}^U, \underline{\mathbf{x}}^W | \underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}, m, r, l) f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m, r, l) \\ &= f_{\text{indep}}(\underline{\mathbf{x}}^W | l) f_{\text{reg}}(\underline{\mathbf{x}}^U | r, \underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j) f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m, r, l) \end{aligned}$$

where $f_{\text{indep}}(\underline{\mathbf{x}}^W | l)$ is the integrated likelihood on the independent irrelevant variable subset W and $f_{\text{reg}}(\underline{\mathbf{x}}^U | r, \underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j)$ corresponds to the integrated likelihood on the subset U regressed on variable subset S and the candidate variable j . The expression of the integrated likelihood restricted on $S \cup \{j\}$ depends on the three situations which can occur for the candidate variable j :

- *First situation:* Given $\underline{\mathbf{x}}^S$, $\underline{\mathbf{x}}^j$ provides additional information for the discriminant analysis thus

$$f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m, r, l) = f_{\text{da}}(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m)$$

corresponds to the integrated likelihood for the discriminant analysis on variable subset $S \cup \{j\}$.

- *Second situation:* Given $\underline{\mathbf{x}}^S$, $\underline{\mathbf{x}}^j$ does not provide additional information for the discriminant analysis but has a linear link with a nonempty subset denoted $R[j]$ of S containing the relevant variables for the regression of $\underline{\mathbf{x}}^j$ on $\underline{\mathbf{x}}^S$:

$$f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m, r, l) = f_{\text{da}}(\underline{\mathbf{x}}^S, \underline{z}|m) f_{\text{reg}}(\underline{\mathbf{x}}^j | [LI], \underline{\mathbf{x}}^{R[j]}).$$

The second term in the right-hand side corresponds to the integrated likelihood of the regression of $\underline{\mathbf{x}}^j$ on $\underline{\mathbf{x}}^{R[j]}$. Since j is a single variable, the variance matrix is spherical ($[LI]$).

- *Third situation:* Given $\underline{\mathbf{x}}^S$, $\underline{\mathbf{x}}^j$ does not provide additional information for the discriminant analysis and is independent of all the variables of S :

$$f(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m, r, l) = f_{\text{da}}(\underline{\mathbf{x}}^S | m) f_{\text{indep}}(\underline{\mathbf{x}}^j | [LI]).$$

The second term in the right-hand side corresponds to the integrated likelihood of the independent Gaussian density on the variable j with a spherical variance matrix since j is a single variable.

In order to compare those three situations in an efficient way, we remark that $f_{\text{indep}}(\underline{\mathbf{x}}^j | [LI])$ can be written $f_{\text{reg}}(\underline{\mathbf{x}}^j | [LI], \underline{\mathbf{x}}^\emptyset)$. Thus instead of considering the nonempty subset $R[j]$ we consider a new explicative variable subset denoted $\tilde{R}[j]$ and defined by $\tilde{R}[j] = \emptyset$ if j follows the third situation and $\tilde{R}[j] = R[j]$ if j follows the second situation. This allows us to recast the comparison of the three situations into the comparison of two situations with the Bayes factor

$$\frac{f_{\text{da}}(\underline{\mathbf{x}}^S, \underline{\mathbf{x}}^j, \underline{z}|m)}{f_{\text{da}}(\underline{\mathbf{x}}^S, \underline{z}|m) f_{\text{reg}}(\underline{\mathbf{x}}^j | [LI], \underline{\mathbf{x}}^{\tilde{R}[j]})}.$$

This Bayes factor being difficult to evaluate, it is approximated by

$$\text{BIC}_{\text{diff}}(j|m) = \text{BIC}_{\text{da}}(\mathbf{x}^S, \mathbf{x}^j, \mathbf{z}|m) - \left\{ \text{BIC}_{\text{da}}(\mathbf{x}^S, \mathbf{z}|m) + \text{BIC}_{\text{reg}}(\mathbf{x}^j | [LI], \mathbf{x}^{\hat{R}[j]}) \right\}. \quad (5)$$

6.2 The general steps of the algorithm

First, this algorithm consists of separating variables into relevant and irrelevant variables for the discriminant analysis via a forward stepwise algorithm described in Section 6.3. Second, the irrelevant variables are partitioned into redundant variables, if regressors are chosen inside the relevant variables, and independent variables otherwise. It remains then to determine the set of regressors from the relevant variables for the multidimensional regression of the redundant variables and the general variance structures.

► For each mixture form m :

- The variable partition into relevant and irrelevant variables for the discriminant analysis, $\hat{S}(m)$ and $\hat{S}^c(m)$ respectively, is determined by the forward stepwise selection algorithm described hereafter (see Section 6.3).
- The variable subset $\hat{S}^c(m)$ is divided into $\hat{U}(m)$ and $\hat{W}(m)$: For each variable j belonging to $\hat{S}^c(m)$, the variable subset $\hat{R}[j]$ of $\hat{S}(m)$ allowing to explain j by a linear regression is determined with the forward stepwise regression algorithm (see Appendix D). If $\hat{R}[j] = \emptyset$, $j \in \hat{W}(m)$ and otherwise, $j \in \hat{U}(m)$.
- For each form r :
 - * The variable subset $\hat{R}(m, r)$, included into $\hat{S}(m)$ and explaining the variables of $\hat{U}(m)$, is determined using the forward stepwise regression algorithm with the fixed form regression model r (see Appendix D).
 - * For each form l : $\hat{\theta}$ and the following criterion value are computed

$$\widetilde{\text{crit}}(m, r, l) = \text{crit}(m, r, l, \hat{S}(m), \hat{R}(m, r), \hat{U}(m), \hat{W}(m)).$$

► The model satisfying the following condition is then selected

$$(\hat{m}, \hat{r}, \hat{l}) = \underset{(m, r, l) \in \mathcal{M} \times \mathcal{T}_{\text{reg}} \times \mathcal{T}_{\text{indep}}}{\text{argmax}} \quad \widetilde{\text{crit}}(m, r, l).$$

► Finally, the complete selected model is

$$\left(\hat{m}, \hat{r}, \hat{l}, \hat{S}(\hat{m}), \hat{R}(\hat{m}, \hat{r}), \hat{U}(\hat{m}), \hat{W}(\hat{m}) \right).$$

Our variable selection procedure is based on forward stepwise algorithms which allow to study data sets where the individual number n is smaller than the variable number Q . Nevertheless, for studying a data set where $Q \leq n$, a backward procedure (starting the search with all variables) could be preferred because it takes variable interactions into account.

6.3 The forward stepwise selection algorithm

Initialisation Let m fixed, $S(m) = \emptyset$, $S^c(m) = \{1, \dots, Q\}$, $j_I = \emptyset$ and $j_E = \emptyset$. The algorithm is making use of an inclusion and an exclusion steps now described. The decision of including (resp. excluding) a variable in (resp. from) the discriminant variable subset is based on (5). Starting from a preliminary inclusion step, the forward variable selection algorithm consists of alternating inclusion and exclusion steps. It returns the discriminant variable subset $\hat{S}(m)$ and the irrelevant variable subset $\hat{S}^c(m)$. These different steps are now described.

Preliminary inclusion step This step consists of selecting the first discriminant variable. For all j in $S^c(m)$, compute

$$\text{BIC}_{\text{diff}}(j|m) = \text{BIC}_{\text{da}}(\mathbf{x}^j, \mathbf{z}|m) - \text{BIC}_{\text{reg}}(\mathbf{x}^j|[LI], \mathbf{x}^\emptyset)$$

and determine

$$j_I = \underset{j \in S^c(m)}{\text{argmax}} \text{BIC}_{\text{diff}}(j|m).$$

Then $S(m) = \{j_I\}$, $S^c(m) = S^c(m) \setminus \{j_I\}$ and go to the inclusion step.

Inclusion step For all j in $S^c(m)$, use the forward stepwise regression algorithm (see Appendix D) to determine the subset $\tilde{R}[j]$ for the regression of \mathbf{x}^j on $\mathbf{x}^{S(m)}$. And, compute

$$\text{BIC}_{\text{diff}}(j|m) = \text{BIC}_{\text{da}}(\mathbf{x}^{S(m)}, \mathbf{x}^j, \mathbf{z}|m) - \left\{ \text{BIC}_{\text{da}}(\mathbf{x}^S, \mathbf{z}|m) + \text{BIC}_{\text{reg}}(\mathbf{x}^j|[LI], \mathbf{x}^{\tilde{R}[j]}) \right\}.$$

Then, compute

$$j_I = \underset{j \in S^c(m)}{\text{argmax}} \text{BIC}_{\text{diff}}(j|m).$$

- If $\text{BIC}_{\text{diff}}(j_I|m) > 0$, $S(m) = S(m) \cup \{j_I\}$, $S^c(m) = S^c(m) \setminus \{j_I\}$ and, if $j_I \neq j_E$, go to the exclusion step and stop otherwise.
- Otherwise, $j_I = \emptyset$. If $j_E \neq \emptyset$, go to the exclusion step and stop otherwise.

Exclusion step For all j in $S(m)$, use the forward stepwise regression algorithm (see Appendix D) to determine the subset $\tilde{R}[j]$ for the regression of \mathbf{x}^j on $\mathbf{x}^{S(m) \setminus \{j\}}$. And, compute

$$\text{BIC}_{\text{diff}}(j|m) = \text{BIC}_{\text{da}}(\mathbf{x}^{S(m)}, \mathbf{z}|m) - \left\{ \text{BIC}_{\text{da}}(\mathbf{x}^{S(m) \setminus \{j\}}, \mathbf{z}|m) + \text{BIC}_{\text{reg}}(\mathbf{x}^j|[LI], \mathbf{x}^{\tilde{R}[j]}) \right\}.$$

Then, compute

$$j_E = \underset{j \in S(m)}{\text{argmin}} \text{BIC}_{\text{diff}}(j|m).$$

- If $\text{BIC}_{\text{diff}}(j_I|m) < 0$, $S^c(m) = S^c(m) \cup \{j_E\}$, $S(m) = S(m) \setminus \{j_E\}$. If $j_E \neq j_I$, go to the inclusion step and stop otherwise.
- Otherwise, $j_E = \emptyset$. If $j_I \neq \emptyset$, go to the inclusion step and stop otherwise.

7 Applications

We present numerical experiments to assess our variable selection procedure. First, the interest of variable selection for non linear discriminant analysis models such as QDA is highlighted on simulated data. Then, two applications on real data sets are presented. The application on the Landsat Satellite data set allows us to illustrate the interest of precisising the role of the variables in an explicative perspective and again the great interest of our variable selection procedure to improve the classification performances of QDA. The second application concerns the Leukemia data of Golub et al. (1999), a classical genomics example where the number of variables is greater than the number of observations.

7.1 Simulated example

This simulated example consists of considering samples described by $Q = 16$ variables. The prior probabilities of the four classes are assumed to be $p_1 = 0.15$, $p_2 = 0.3$, $p_3 = 0.2$ and $p_4 = 0.35$. On the three discriminant variables, data are distributed from $\mathbf{x}_i^{\{1-3\}}|z_i = k \sim \Phi(\cdot|\mu_k, \Sigma_k)$ with $\mu_1 = (1.5, -1.5, 1.5)$, $\mu_2 = (-1.5, 1.5, 1.5)$, $\mu_3 = (1.5, -1.5, -1.5)$, $\mu_4 = (-1.5, 1.5, -1.5)$, and $\Sigma_k = (\rho_k^{|i-j|})_{1 \leq i,j \leq 3}$ with $\rho_1 = 0.85$, $\rho_2 = 0.1$, $\rho_3 = 0.65$ and $\rho_4 = 0.5$. Four redundant variables simulated from

$$\mathbf{x}_i^{\{4-7\}} \sim \mathcal{N}\left(\mathbf{x}_i^{\{1,3\}} \begin{pmatrix} 1 & 0 & -1 & 2 \\ 0 & -2 & 2 & 1 \end{pmatrix}; I_4\right)$$

and nine independent variables are appended, sampled from $\mathbf{x}_i^{\{8-16\}} \sim \mathcal{N}(\gamma, \tau)$ with

$$\gamma = (-2, -1.5, -1, -0.5, 0, 0.5, 1, 1.5, 2)$$

and

$$\tau = \text{diag}(0.5, 0.75, 1, 1.25, 1.5, 1.25, 1, 0.75, 0.5).$$

A total of 100 simulation replications are considered where the training sample is composed of $n = 500$ data points and the same test sample with 50,000 points is used. The LDA, QDA and EDDA methods with and without variable selection are compared according to the averaged classification error rates.

Results summarized in Table 1 show that the variable selection procedure allows to improve the classification performance of LDA, QDA and EDDA. In particular, QDA becomes superior to LDA with variable selection. In all replications, only the first three variables are declared discriminant. When the true variable partition is not selected, it is due to some independent variables which are declared redundant (24, 24 and 26 times for EDDA, QDA and LDA respectively).

| with variable selection | | | without variable selection | | |
|-------------------------|------------|------------|----------------------------|------------|------------|
| LDA | QDA | EDDA | LDA | QDA | EDDA |
| 4.94 | 4.19 | 4.18 | 5.30 | 6.23 | 5.29 |
| ± 0.13 | ± 0.06 | ± 0.06 | ± 0.18 | ± 0.38 | ± 0.18 |

Table 1: Averaged classification error rate (\pm standard deviation) for LDA, QDA and EDDA methods, with and without variable selection for the simulated data sets.

7.2 The Landsat Satellite Data

The Landsat Satellite Data, available at the UCI Machine Learning Repository (see <http://www.ics.uci.edu/~mllearn/>) is considered. This data set consists of the multi-spectral values of pixels in a tiny sub-area of a satellite image. Each line is a vector of length $Q = 36$, composed of the pixel values in four spectral bands (two in the visible region and two in the near infra-red) of each of the 9 pixels in the 3×3 neighborhood. These data points are split into six classes. The original data set has already been divided into a training set with 4,435 samples and a testing set with 2,000 samples. The same experiment conditions than in Zhang and Wang (2008) are considered: 1,000 samples (randomly selected from the training data) are used to estimate and select the model, and this experiment is randomly replicated 100 times. Only QDA and LDA are considered in this study.

According to Table 2, QDA and LDA perform the same without variable selection, while QDA outperforms LDA with variable selection. In all replications, our variable selection procedure selects the QDA model ($\hat{m} = [L_k C_k]$), and all the irrelevant classification variables are redundant ($\hat{W} = \emptyset$) and regressed on all discriminant variables ($\hat{R} = \hat{S}$) with a general covariance matrix structure ($\hat{r} = [LC]$).

| with variable selection | | without variable selection | |
|-------------------------|------------|----------------------------|------------|
| LDA | QDA | LDA | QDA |
| 21.00 | 16.21 | 18.05 | 17.90 |
| ± 0.53 | ± 0.68 | ± 0.48 | ± 0.57 |

Table 2: Averaged classification error rate (\pm standard deviation) for LDA and QDA methods, without and with variable selection for the Landsat Satellite Data.

It is noteworthy that QDA and LDA select the same variables in the 100 replications, with an average selection of 12 discriminant variables as in Zhang and Wang (2008). It is worth mentioning some variable selection tendencies (see Figure 1). First, variables tend to be selected by couple: for instance Variables 34 and 36, Variables 18 and 20, and Variables 2 and 4 are both declared discriminant in the same replications. Second, we can note that Variables 3, 7, 11, 15, 19, 23, 27, 31 and 35, corresponding to one measure in the near infra-red for each pixel of the 3×3 neighborhood are never declared discriminant. Third, the variables corresponding to Pixels 1, 3, 5, 7 and 9 in the 3×3 neighborhood are more often declared discriminant than the one of the other pixels, certainly because these pixels have more neighbor pixels in common.

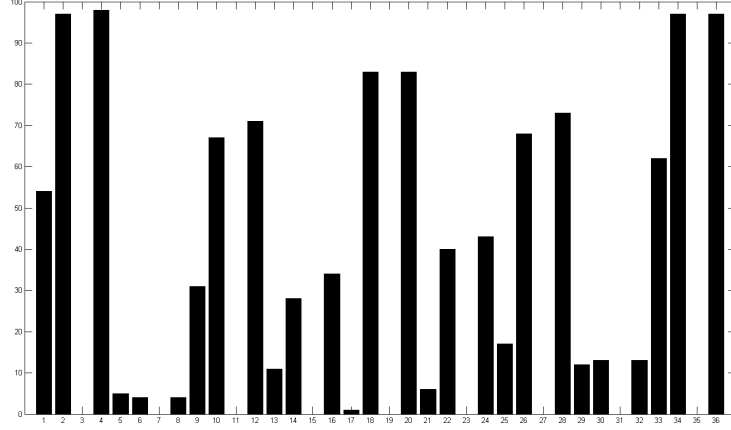


Figure 1: Number of times each variable is declared discriminant among 100 replications for the Landsat Satellite Data.

7.3 Leukemia data set

These data come from a study of gene expression in two types of acute leukemias: acute lymphoblastic leukemia (ALL) and acute myeloid leukemia (AML) published by Golub et al. (1999). Gene expression levels were measured for 47 ALL tumor samples and 25 AML tumor samples, using Affymetrix high-density oligonucleotide arrays containing 6817 human genes. After the pre-processing steps, as image analysis, standardization and some gene filtering, $Q = 3571$ genes are conserved. The interest of this data set is the large number of genes describing the samples, and the importance to detect the genes whose expression pattern separate the two types of leukemias. This data set is known as a benchmark data set and numerous results are available (Golub et al., 1999; Su et al., 2003; Krishnapuram et al., 2004; Mary-Huard and Robin, 2009; Yang and Xin-Yuan, 2010). It is considered from two points of view in this section. First, the relevance and stability of our variable selection procedure are measured using a leave-one-out procedure for LDA and QDA models. Second, the interest of variable selection to improve the prediction accuracy is assessed using the training and test samples used by Golub et al. (1999).

In the leave-one-out procedure for LDA and QDA, the averaged classification error rate is equal to zero and so as good as the error rates of different methods already applied (see Mary-Huard and Robin, 2009). For LDA, 3529 genes are never declared discriminant for the classification by our variable selection procedure and 44 of them are always declared independent. Among the 42 genes declared at least once discriminant, seven genes (Macmarcks, CD33 antigen, CST3, DF D, CCND3, GLUTATHIONE S-TRANSFERASE and RABAPTIN-5 protein) are already known to be discriminant and three (PLZF, Adrenal-Specific Protein Pg2 and, PRSS1) are known to be implicated in cancers. The first two genes declared the most time relevant are DF D (67 times) and GLUTATHIONE S-TRANSFERASE (53 times); the other ones are more declared redundant than relevant. We also compare our results with the list of Su et al.

(2003) which contains 100 genes reported as discriminant by at least one discriminant method. A lack of precision in the gene names leads to consider 88 genes among these 100 genes. According to our method, 83 genes are mainly declared redundant and six of them (Macmarcks, CD33 antigen, CST3, DF D (adipsin), CCND3 and GLUTATHIONE S-TRANSFERASE) at least once relevant. The five remaining genes are declared at least once independent (W) and never discriminant for the classification. Among these five genes, three are very often in W (more than 52 times) but are identified in the Su list by only the t-test procedure. They may be false candidats. The two others (Lamp2 and GLYCYLPEPTIDE TETRADECANOYLTRANSFERASE) are declared redundant 67 and 60 times and otherwise independent. Their status is thus more redundant than independent and it is coherent with their presence in the Su list. We do not report in detail the results in QDA which are quite analogous. But it is worthwhile to remark that the number of discriminant variables selected at least one time is 122 for QDA instead of 32 for LDA. It indicates a greater stability of the variable selection procedure for the simpler model LDA.

We then analyze the Leukemia data set using 38 (27 are ALL and 11 are AML) samples in the training and 34 (20 are ALL and 14 are AML) samples in the test. Results of our method are given in Table 3 and are compared with the performance of other methods given in Yang and Xin-Yuan (2010). Despite the variable selection, LDA performs poorly, on the contrary the quadratic methods (QDA and $[L_k C]$) are greatly improved by variable selection leading to zero misclassified test observation. Moreover, this is achieved with a small number of discriminant variables especially for the parsimonious model $[L_k C]$ (see Appendix A).

| Models | LDA | QDA | $[L_k C]$ |
|----------------------------|-------|-------|-----------|
| $\text{card}(\hat{S})$ | 8 | 8 | 3 |
| $\text{card}(\hat{R})$ | 2 | 2 | 3 |
| $\text{card}(\hat{U})$ | 3058 | 2848 | 1912 |
| $\text{card}(\hat{W})$ | 505 | 715 | 1656 |
| Misc. test obs. (ALL, AML) | (2,4) | (0,0) | (0,0) |

Table 3: Variable selection and misclassification error rate. The first four lines indicate the number of variables in S , R , U and W sets. The last line gives the number of misclassified test observations according to the two leukemia types ALL and AML.

8 Discussion

We have proposed a variable selection methodology for a large family of Gaussian generative models in discriminant analysis. Regarding the problem as a model selection problem, we have proposed a BIC-like criterion to distinguish between discriminant, redundant and noisy variables. We proved the identifiability of our model collection and the consistency of the proposed BIC-like criterion. A procedure embedding two forward stepwise variable selection algorithms for classification and regression has been defined. Numerical experiments highlight the potentially great interest of our variable selection procedure to improve the classification performances of non linear Gaussian classification

models. Actually those models involve many parameters when the number of variables is large with respect to the training sample size. But our variable selection procedure allows us to overcome the dimensionality problem leading to powerful classifiers with a nice interpretation of variable roles. Those results confirm the promising performances obtained by Murphy et al. (2010) and Zhang and Wang (2008) with less general methods. Our opinion is that our methodology is able to make the non linear generative classification methods such as quadratic discriminant analysis much more efficient in high dimensional contexts and competitive with gold standard classifiers such as LDA, logistic regression, k -nearest neighbor classifier or support vector classifiers in many situations.

Appendices

A The different model forms

This is the list of the 14 different model forms available in the MIXMOD software.

| Family | Model | Volume | Orientation | Shape |
|-----------|------------------|----------|-----------------|----------|
| Spherical | $[LI]$ | equal | NA | equal |
| | $[L_k I]$ | variable | NA | equal |
| Diagonal | $[LB]$ | equal | coordinate axes | equal |
| | $[L_k B]$ | variable | coordinate axes | equal |
| | $[LB_k]$ | equal | coordinate axes | variable |
| | $[L_k B_k]$ | variable | coordinate axes | variable |
| General | $[LC]$ | equal | equal | equal |
| | $[L_k C]$ | variable | equal | equal |
| | $[LDA_k D]$ | equal | equal | variable |
| | $[L_k DA_k D]$ | variable | equal | variable |
| | $[LD_k AD_k]$ | equal | variable | equal |
| | $[L_k D_k AD_k]$ | variable | variable | equal |
| | $[LC_k]$ | equal | variable | variable |
| | $[L_k C_k]$ | variable | variable | variable |

Table 4: List of model forms available in MIXMOD.

B Proof of the model identifiability

Theorem 1 concerning the model identifiability can be proved quickly (see Proof 1) using the SRUW model identifiability established in Maugis et al. (2009b) in the model-based clustering context. It is also possible to completely prove this theorem in the discriminant analysis context (see Proof 2).

Proof 1. Let (m, r, l, \mathbf{V}) and $(m^*, r^*, l^*, \mathbf{V}^*)$ be two models. Let $\theta \in \Theta_{(m, r, l, \mathbf{V})}$ and $\theta^* \in \Theta_{(m^*, r^*, l^*, \mathbf{V}^*)}$ two parameters vectors such that $\forall \mathbf{x} \in \mathbb{R}^Q, \forall z \in \{1, \dots, K\}$,

$$f(\mathbf{x}, z | m, r, l, \mathbf{V}, \theta) = f(\mathbf{x}, z | m^*, r^*, l^*, \mathbf{V}^*, \theta^*).$$

It is equivalent to the following system: $\forall \mathbf{x} \in \mathbb{R}^Q, \forall k \in \{1, \dots, K\}$,

$$\begin{aligned} & p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)}) \Phi(\mathbf{x}^W | \gamma, \tau_{(l)}) \\ = & p_k^* \Phi(\mathbf{x}^S | \mu_k^*, \Sigma_{k(m^*)}^*) \Phi(\mathbf{x}^U | a^* + \mathbf{x}^R \beta^*, \Omega_{(r^*)}^*) \Phi(\mathbf{x}^W | \gamma^*, \tau_{(l^*)}^*) \end{aligned} \quad (6)$$

Summing the K previous equations, we obtain that

$$\begin{aligned} & \left\{ \sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \right\} \Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)}) \Phi(\mathbf{x}^W | \gamma, \tau_{(l)}) \\ &= \left\{ \sum_{k=1}^K p_k^* \Phi(\mathbf{x}^S | \mu_k^*, \Sigma_{k(m^*)}^*) \right\} \Phi(\mathbf{x}^U | a^* + \mathbf{x}^R \beta^*, \Omega_{(r^*)}^*) \Phi(\mathbf{x}^W | \gamma^*, \tau_{(l^*)}^*). \end{aligned}$$

Next, using the identifiability result established in Maugis et al. (2009b) in the clustering framework with a fix number of components K , we obtain that $m = m^*$, $r = r^*$, $l = l^*$, $a = a^*$, $\beta = \beta^*$, $\Omega_{(r)} = \Omega_{(r^*)}^*$ and the parameters p_k , p_k^* , μ_k , μ_k^* and $\Sigma_{k(m)}^*$, $\Sigma_{k(m^*)}^*$ are equal up to a permutation of Gaussian mixture components. But this permutation is the identity according to (6) thus $p_k = p_k^*$, $\mu_k = \mu_k^*$ and $\Sigma_{k(m)}^* = \Sigma_{k(m^*)}^*$ for all $k \in \{1, \dots, K\}$.

Proof 2. Let (m, r, l, \mathbf{V}) and $(m^*, r^*, l^*, \mathbf{V}^*)$ be two models. Let $\theta \in \Theta_{(m, r, l, \mathbf{V})}$ and $\theta^* \in \Theta_{(m^*, r^*, l^*, \mathbf{V}^*)}$ two parameters vectors such that $\forall \mathbf{x} \in \mathbb{R}^Q$, $\forall z \in \{1, \dots, K\}$,

$$f(\mathbf{x}, z | m, r, l, \mathbf{V}, \theta) = f(\mathbf{x}, z | m^*, r^*, l^*, \mathbf{V}^*, \theta^*).$$

It is equivalent to the following system: $\forall \mathbf{x} \in \mathbb{R}^Q$, $\forall k \in \{1, \dots, K\}$,

$$\begin{aligned} & p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)}) \Phi(\mathbf{x}^W | \gamma, \tau_{(l)}) \\ &= p_k^* \Phi(\mathbf{x}^S | \mu_k^*, \Sigma_{k(m^*)}^*) \Phi(\mathbf{x}^U | a^* + \mathbf{x}^R \beta^*, \Omega_{(r^*)}^*) \Phi(\mathbf{x}^W | \gamma^*, \tau_{(l^*)}^*) \end{aligned}$$

and can be reformulated as

$$p_k \Phi(\mathbf{x} | A_k, B_k) = p_k^* \Phi(\mathbf{x} | A_k^*, B_k^*) \quad (7)$$

where the Q -dimensional vectors A_k are defined by

$$\forall j \in \{1, \dots, Q\}, \quad A_{kj} = \begin{cases} \mu_{kj} & \text{if } j \in S \\ (d + \mu_k \Lambda)_j & \text{if } j \in S^c \end{cases}$$

and the $Q \times Q$ -dimensional matrices B_k are defined by

$$\forall (i, j) \in \{1, \dots, Q\}^2, \quad B_{k,ij} = \begin{cases} \Sigma_{k(m),ij} & \text{if } i \in S, j \in S \\ (\Sigma_{k(m)} \Lambda)_{ij} & \text{if } i \in S, j \in S^c \\ (\Lambda' \Sigma_{k(m)})_{ij} & \text{if } i \in S^c, j \in S \\ (D + \Lambda' \Sigma_{k(m)} \Lambda)_{ij} & \text{if } i \in S^c, j \in S^c \end{cases}$$

with

$$\forall j \in S^c, \quad d_j = \begin{cases} a_j & \text{if } j \in U \\ \gamma_j & \text{if } j \in W \end{cases}$$

$$\forall j \in S^c, \quad \forall h \in S, \quad \Lambda_{jh} = \begin{cases} \beta_{jh} & \text{if } j \in U, h \in R \\ 0 & \text{otherwise} \end{cases}$$

and

$$\forall (i, j) \in (S^c)^2, \quad D_{ij} = \begin{cases} \Omega_{(r),ij} & \text{if } i \in U, j \in U \\ \tau_{(l),ij} & \text{if } i \in W, j \in W \\ 0 & \text{otherwise.} \end{cases}$$

In the same way, we define the A_k^* 's and B_k^* 's. In order to make easier the reading of this proof, the indexation of Σ_k , Ω and τ by m , r and l are omitted.

First according to (7), for all $k \in \{1, \dots, K\}$, $p_k = p_k^*$, $A_k = A_k^*$ and $B_k = B_k^*$. Indeed if there exists k such that $p_k < p_k^*$ then $\forall \mathbf{x} \in \mathbb{R}^Q$, $\Phi(\mathbf{x}|A_k, B_k) > \Phi(\mathbf{x}|A_k^*, B_k^*)$ and that is in contradiction with the fact that $\Phi(\cdot|A_k, B_k)$ and $\Phi(\cdot|A_k^*, B_k^*)$ are two densities.

Second, assume that $S \cap S^* = \emptyset$ and consider the subsets $s = S^* \cap S^c$ and $t = S^{*c} \cap S$. The equality of variance matrices B_k and B_k^* on s and between s and t gives respectively for all $k \in \{1, \dots, K\}$,

$$\begin{cases} D_{ss} + \Lambda'_{.s} \Sigma_k \Lambda_{.s} = \Sigma_k^* \\ \Lambda'_{.s} \Sigma_k = \Sigma_k^* \Lambda_{.t}^* \end{cases}.$$

According to these two equalities, it is deduced that $D_{ss} = \Sigma_k^*(I - \Lambda_{.t}^* \Lambda_{.s})$ for all k . Since D_{ss} and Σ_k^* are positive definite matrices, $I - \Lambda_{.t}^* \Lambda_{.s}$ is a nonsingular matrix and consequently all the variance matrices $\Sigma_k^* = D_{ss}(I - \Lambda_{.t}^* \Lambda_{.s})^{-1}$ are equal. Moreover, the equality of mean vectors on s and t gives

$$\begin{cases} d_s + \mu_k \Lambda_{.s} = \mu_k^* \\ \mu_k = d_t^* + \mu_k^* \Lambda_{.t}^* \end{cases},$$

implying that $\mu_k^*(I - \Lambda_{.t}^* \Lambda_{.s}) = d_s + d_t^* \Lambda_{.s}$, for all k . Since $I - \Lambda_{.t}^* \Lambda_{.s}$ is nonsingular, all μ_k^* are equal. This is in contradiction with the assumption (C1).

Third, assume that $S \cap S^* \neq \emptyset$ and $S^c \cap S^{*c} \neq \emptyset$, and consider the nonempty subsets $t = S^c \cap S^*$, $s = S \cap S^*$ and $\bar{s} = S \cap S^{*c}$. The equality of variance matrices B_k and B_k^* on \bar{s} , on s , between t and s , between s and \bar{s} , and between t and \bar{s} gives respectively for all k ,

$$\begin{aligned} \Sigma_{k,\bar{s}\bar{s}} &= D_{\bar{s}\bar{s}}^* + \Lambda_{\bar{s}\bar{s}}^{*'}(\Sigma_{k,ss}^* \Lambda_{\bar{s}\bar{s}}^* + \Sigma_{k,st}^* \Lambda_{t\bar{s}}^*) + \Lambda_{\bar{s}\bar{s}}^{*'}(\Sigma_{k,ts}^* \Lambda_{\bar{s}\bar{s}}^* + \Sigma_{k,tt}^* \Lambda_{t\bar{s}}^*) \\ \Sigma_{k,ss} &= \Sigma_{k,ss}^* \end{aligned} \quad (9)$$

$$\Lambda_{\bar{s}t}' \Sigma_{k,\bar{s}s} + \Lambda_{st}' \Sigma_{k,ss} = \Sigma_{k,ts}^* \quad (10)$$

$$\Sigma_{k,\bar{s}s} = \Lambda_{\bar{s}\bar{s}}^{*'} \Sigma_{k,ss}^* + \Lambda_{\bar{s}\bar{s}}^{*'} \Sigma_{k,ts}^* \quad (11)$$

$$\Sigma_{k,\bar{s}\bar{s}} \Lambda_{\bar{s}t} + \Sigma_{k,\bar{s}s} \Lambda_{st} = \Lambda_{\bar{s}\bar{s}}^{*'} \Sigma_{k,st}^* + \Lambda_{\bar{s}\bar{s}}^{*'} \Sigma_{k,tt}^* \quad (12)$$

From (8), (11), (12), we get

$$\Sigma_{k,\bar{s}\bar{s}}(I - \Lambda_{\bar{s}t} \Lambda_{t\bar{s}}^*) = D_{\bar{s}\bar{s}}^* + \Sigma_{k,\bar{s}s}(\Lambda_{\bar{s}\bar{s}}^* + \Lambda_{st}^* \Lambda_{t\bar{s}}^*) \quad (13)$$

and Equations (9), (10) and (11) allow to deduce

$$\Lambda_{\bar{s}\bar{s}}^* + \Lambda_{st}^* \Lambda_{t\bar{s}}^* = \Sigma_{k,ss}^{-1} \Sigma_{k,\bar{s}s}(I - \Lambda_{\bar{s}t} \Lambda_{t\bar{s}}^*). \quad (14)$$

Finally, Equations (13) and (14) imply $\Sigma_{k,\bar{s}\bar{s}}(I - \Lambda_{\bar{s}t} \Lambda_{t\bar{s}}^*) = D_{\bar{s}\bar{s}}^*$. Since $D_{\bar{s}\bar{s}}^*$ and $\Sigma_{k,\bar{s}\bar{s}}|_s$ are positive definite matrices, the matrix $I - \Lambda_{\bar{s}t} \Lambda_{t\bar{s}}^*$ is nonsingular and all the matrices $\Sigma_{k,\bar{s}\bar{s}}|_s$ are equal. Similarly, according to (14), all matrices $\Sigma_{k,\bar{s}}|_s$ are equal. The equality of mean vectors on \bar{s} , s and t gives the following equations: For all k ,

$$\begin{cases} \mu_{k\bar{s}} = d_{\bar{s}}^* + \mu_{ks}^* \Lambda_{\bar{s}\bar{s}}^* + \mu_{k,t}^* \Lambda_{t\bar{s}}^* \\ \mu_{ks} = \mu_{ks}^* \\ d_t + \mu_{k\bar{s}} \Lambda_{\bar{s}t} + \mu_{ks} \Lambda_{st} = \mu_{k,t}^* \end{cases}$$

implying

$$\mu_{k\bar{s}}(I - \Lambda_{\bar{s}t} \Lambda_{t\bar{s}}^*) = (d_{\bar{s}}^* + d_t \Lambda_{t\bar{s}}^*) + \mu_{ks}(\Lambda_{\bar{s}\bar{s}}^* + \Lambda_{st}^* \Lambda_{t\bar{s}}^*)$$

and Equation (14) leads to

$$(\mu_{k\bar{s}} - \mu_{ks} \Sigma_{k,ss}^{-1} \Sigma_{k,s\bar{s}})(I - \Lambda_{st} \Lambda_{t\bar{s}}^*) = d_{\bar{s}} + d_t^* \Lambda_{t\bar{s}}^*.$$

Since $I - \Lambda_{st} \Lambda_{t\bar{s}}^*$ is non singular, the mean vectors $\mu_{k,\bar{s}|s}$ are also equal, and thus the constraint (C1) is violated. In the same way, assuming that \bar{s} or t is empty, we prove that $S \subsetneq S^*$ and $S^* \subsetneq S$ are impossible.

Finally, it leads to $S = S^*$ and, by the equality of variance matrices and mean vectors, we easily obtain that $\mu_k = \mu_k^*$, $\Sigma_k = \Sigma_k^*$ (and then $m = m^*$), $d = d^*$, $D = D^*$ and $\Lambda = \Lambda^*$. Then $R = R^*$ because otherwise, according to the definition of Λ , there exists $h \in R \cap R^{*c}$ such that $\forall u \in U, \beta_{hu} = 0$ or there exists $h \in R^* \cap R^c$ such that $\forall u \in U^*, \beta_{hu}^* = 0$ that is contradicted the assumption (C2). In same way, we prove that $U = U^*$ and thus $W = W^*$. Finally, according to the definition of Λ , d and D , we obtain that $\beta = \beta^*$, $\Omega = \Omega^*$, $r = r^*$, $a = a^*$, $\gamma = \gamma^*$, $\tau = \tau^*$ and $l = l^*$.

C Proof of the criterion consistency theorem

This appendix is devoted to the proof of Theorem 2 given the criterion consistency.

Proof. According to the expressions (2) and (3), the selected model satisfies

$$(\hat{m}, \hat{r}, \hat{l}, \hat{\mathbf{V}}) = \underset{(m,r,l,\mathbf{V}) \in \mathcal{N}}{\operatorname{argmax}} \operatorname{crit}(m, r, l, \mathbf{V})$$

with

$$\operatorname{crit}(m, r, l, \mathbf{V}) = 2 \sum_{i=1}^n \ln[f(\mathbf{x}_i, z_i | m, r, l, \mathbf{V}, \hat{\theta})] - \Xi_{(m,r,l,\mathbf{V})} \ln(n).$$

Thus

$$\begin{aligned} P((\hat{m}, \hat{r}, \hat{l}, \hat{\mathbf{V}}) = (m_0, r_0, l_0, \mathbf{V}_0)) \\ = P(\operatorname{crit}(m_0, r_0, l_0, \mathbf{V}_0) - \operatorname{crit}(m, r, l, \mathbf{V}) \geq 0, \forall (m, r, l, \mathbf{V}) \in \mathcal{N}). \end{aligned} \quad (15)$$

Denoting $\Delta \operatorname{crit}(m, r, l, \mathbf{V}) = \operatorname{crit}(m_0, r_0, l_0, \mathbf{V}_0) - \operatorname{crit}(m, r, l, \mathbf{V})$, we get

$$\begin{aligned} \Delta \operatorname{crit}(m, r, l, \mathbf{V}) &= 2n \left[\frac{1}{n} \sum_{i=1}^n \ln \left\{ \frac{f(\mathbf{x}_i, z_i | m_0, r_0, l_0, \mathbf{V}_0, \hat{\theta})}{h(\mathbf{x}_i, z_i)} \right\} - \frac{1}{n} \sum_{i=1}^n \ln \left\{ \frac{f(\mathbf{x}_i, z_i | m, r, l, \mathbf{V}, \hat{\theta})}{h(\mathbf{x}_i, z_i)} \right\} \right] \\ &\quad + [\Xi_{(m,r,l,\mathbf{V})} - \Xi_{(m_0,r_0,l_0,\mathbf{V}_0)}] \ln(n). \end{aligned} \quad (16)$$

Let $\mathcal{N}_1 = \{(m, r, l, \mathbf{V}) \in \mathcal{N}; \operatorname{KL}[h, f(\cdot | m, r, l, \mathbf{V}, \theta^*)] \neq 0\}$. We have that $\mathcal{N}_1 = \mathcal{N} \setminus \{(m_0, r_0, l_0, \mathbf{V}_0)\}$ since if $\operatorname{KL}[h, f(\cdot | m, r, l, \mathbf{V}, \theta^*)] = 0$ then $h = f(\cdot | m_0, r_0, l_0, \mathbf{V}_0, \theta^*) = f(\cdot | m, r, l, \mathbf{V}, \theta^*)$ and according to the model identifiability, $(m_0, r_0, l_0, \mathbf{V}_0) = (m, r, l, \mathbf{V})$. Thus from (15), the theorem is established if it is proved that

$$\forall (m, r, l, \mathbf{V}) \in \mathcal{N}_1, P(\Delta \operatorname{crit}(m, r, l, \mathbf{V}) < 0) \xrightarrow{n \rightarrow \infty} 0. \quad (17)$$

Let $(m, r, l, \mathbf{V}) \in \mathcal{N}_1$. Denoting $\mathbb{M}_n(m, r, l, \mathbf{V}) = \frac{1}{n} \sum_{i=1}^n \ln \left\{ \frac{f(\mathbf{x}_i, z_i | m, r, l, \mathbf{V}, \hat{\theta})}{h(\mathbf{x}_i, z_i)} \right\}$ and $M(m, r, l, \mathbf{V}) = -\text{KL}[h, f(\cdot | m, r, l, \mathbf{V}, \theta^*)]$, from (16) we have

$$\begin{aligned} & P(\Delta_{\text{crit}}(m, r, l, \mathbf{V}) < 0) \\ &= P(2n\{\mathbb{M}_n(m_0, r_0, l_0, \mathbf{V}_0) - \mathbb{M}_n(m, r, l, \mathbf{V})\} \\ &\quad + [\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n) < 0) \\ &= P(\mathbb{M}_n(m_0, r_0, l_0, \mathbf{V}_0) - M(m_0, r_0, l_0, \mathbf{V}_0) + M(m, r, l, \mathbf{V}) - \mathbb{M}_n(m, r, l, \mathbf{V}) \\ &\quad + M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) + \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} < 0). \end{aligned}$$

Thus, for all $\epsilon > 0$, according to Lemma 7,

$$\begin{aligned} & P(\Delta_{\text{crit}}(m, r, l, \mathbf{V}) < 0) \\ &\leq P(M(m_0, r_0, l_0, \mathbf{V}_0) - \mathbb{M}_n(m_0, r_0, l_0, \mathbf{V}_0) > \epsilon) \\ &\quad + P(\mathbb{M}_n(m, r, l, \mathbf{V}) - M(m, r, l, \mathbf{V}) > \epsilon) \\ &\quad + P\left(M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) + \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} < 2\epsilon\right). \end{aligned}$$

From Lemma 3, stated hereafter, $\forall (m, r, l, \mathbf{V}) \in \mathcal{N}, \mathbb{M}_n(m, r, l, \mathbf{V}) \xrightarrow[n \rightarrow \infty]{P} M(m, r, l, \mathbf{V})$.

Thus,

$$\forall \epsilon > 0, P(\mathbb{M}_n(m, r, l, \mathbf{V}) - M(m, r, l, \mathbf{V}) > \epsilon) \leq P(|\mathbb{M}_n(m, r, l, \mathbf{V}) - M(m, r, l, \mathbf{V})| > \epsilon) \xrightarrow[n \rightarrow \infty]{} 0.$$

For the third term, note

$$\begin{aligned} & P\left(M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) + \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} < 2\epsilon\right) \\ &\leq P\left(M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) - 2\epsilon < \left| \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} \right|\right). \end{aligned}$$

Since $[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)/2n \xrightarrow[n \rightarrow \infty]{} 0$ and $M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) > 0$ because $(m, r, l, \mathbf{V}) \in \mathcal{N}_1$, taking $\epsilon = \{M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V})\}/4 > 0$, we get

$$\begin{aligned} & P\left(M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V}) + \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} < 2\epsilon\right) \\ &\leq P\left(\frac{M(m_0, r_0, l_0, \mathbf{V}_0) - M(m, r, l, \mathbf{V})}{2} < \left| \frac{[\Xi_{(m, r, l, \mathbf{V})} - \Xi_{(m_0, r_0, l_0, \mathbf{V}_0)}] \ln(n)}{2n} \right|\right) \xrightarrow[n \rightarrow \infty]{} 0. \end{aligned}$$

Finally, $P(\Delta_{\text{crit}}(m, r, l, \mathbf{V}) < 0) \xrightarrow[n \rightarrow \infty]{} 0$.

□

Lemma 3. Under assumptions (H1) and (H2),

$$\forall (m, r, l, \mathbf{V}) \in \mathcal{N}, \frac{1}{n} \sum_{i=1}^n \ln \left[\frac{h(\mathbf{x}_i, z_i)}{f(\mathbf{x}_i, z_i | m, r, l, \mathbf{V}, \hat{\theta})} \right] \xrightarrow[n \rightarrow \infty]{P} \text{KL}[h, f(\cdot | m, r, l, \mathbf{V}, \theta^*)].$$

Proof. Let $(m, r, l, \mathbf{V}) \in \mathcal{N}$. By the law of large numbers, if $\mathbb{E}[|\ln(h(X))|] < \infty$,

$$\frac{1}{n} \sum_{i=1}^n \ln[h(\mathbf{x}_i, z_i)] \xrightarrow[n \rightarrow \infty]{P} \mathbb{E}[\ln(h(X, Z))]. \quad (18)$$

And, if the Proposition 4 can be applied with the family

$$\mathcal{F}_{(m,r,l,\mathbf{V})} := \{\ln[f(\cdot|m, r, l, \mathbf{V}, \theta)]; \theta \in \Theta'_{(m,r,l,\mathbf{V})}\}$$

thus

$$\frac{1}{n} \sum_{i=1}^n \ln[f(\mathbf{x}_i, z_i|m, r, l, \mathbf{V}, \hat{\theta})] \xrightarrow[n \rightarrow \infty]{P} \mathbb{E}[\ln f(X, Z|m, r, l, \mathbf{V}, \theta^*)]. \quad (19)$$

Then (18) and (19) give the result. Thus we have to prove that (H2) allows to verify the hypotheses of the Proposition 4 and $\mathbb{E}[|\ln h(X, Z)|] < \infty$.

Firstly, according to (H2), $\Theta'_{(m,r,l,\mathbf{V})}$ is a compact metric space. Moreover, for all \mathbf{x} in $\mathbb{R}^Q, \forall z \in \{1, \dots, K\}$, $\theta \in \Theta'_{(m,r,l,\mathbf{V})} \mapsto \ln[f(\mathbf{x}, z|m, r, l, \mathbf{V}, \theta)]$ is continuous. Let us verify now that there is an envelope function F of $\mathcal{F}_{(m,r,l,\mathbf{V})}$ being h -integrable.

Recalling that

$$\begin{aligned} \ln[f(\mathbf{x}, z|m, r, l, \mathbf{V}, \theta)] &= \ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \mathbb{I}_{z=k} \right] \\ &\quad + \ln[\Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)})] + \ln[\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})], \end{aligned}$$

these three terms are bounded separately.

Study of the first term:

Due to $\|\mathbf{x}^S - \mu_k\|_{\Sigma_k^{-1}}^2 \geq 0$, $|\Sigma_k|^{-\frac{1}{2}} \leq s_m^{-\frac{\#S}{2}}$ according to Lemma 5 and $\sum_{k=1}^K p_k = 1$, the upper bound of this first term is given by

$$\begin{aligned} \ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \mathbb{I}_{z=k} \right] &\leq \ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \right] \\ &\leq \ln \left[\sum_{k=1}^K p_k |2\pi \Sigma_k|^{-\frac{1}{2}} \exp \left(-\frac{\|\mathbf{x}^S - \mu_k\|_{\Sigma_k^{-1}}^2}{2} \right) \right] \\ &\leq \ln \left[\sum_{k=1}^K p_k (2\pi s_m)^{-\frac{\#S}{2}} \right] \\ &\leq -\frac{\#S}{2} \ln [2\pi s_m] \end{aligned}$$

where $\|\mathbf{x}^S - \mu_k\|_{\Sigma_k^{-1}}^2 = (\mathbf{x}^S - \mu_k) \Sigma_k^{-1} (\mathbf{x}^S - \mu_k)'$.

For obtaining a lower bound,

$$\begin{aligned}
\ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \mathbb{I}_{z=k} \right] &= \ln [p_z \Phi(\mathbf{x}^S | \mu_z, \Sigma_{z(m)})] \\
&= \ln(p_z) + \ln \left[|2\pi \Sigma_z|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} \|\mathbf{x}^S - \mu_z\|_{\Sigma_z^{-1}}^2 \right) \right] \\
&= \ln(p_z) - \frac{\#S}{2} \ln[2\pi] - \frac{1}{2} \{ \ln[|\Sigma_z|] + [\|\mathbf{x}^S - \mu_z\|_{\Sigma_z^{-1}}^2] \}.
\end{aligned}$$

Since $|\Sigma_z| \leq s_M^{\#S}$ according to Lemma 5, $p_z \geq \rho$ and

$$\begin{aligned}
\|\mathbf{x}^S - \mu_z\|_{\Sigma_z^{-1}}^2 &\leq \frac{\|\mathbf{x}^S - \mu_z\|^2}{s_m} \\
&\leq \frac{2(\|\mathbf{x}^S\|^2 + \|\mu_z\|^2)}{s_m} \\
&\leq \frac{2(\|\mathbf{x}^S\|^2 + \eta^2)}{s_m}
\end{aligned}$$

because $\mu_z \in \mathcal{B}(\eta, \#S)$, we obtain that

$$\ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \mathbb{I}_{z=k} \right] \geq \ln(\rho) - \frac{\#S}{2} \ln[2\pi s_M] - \frac{\|\mathbf{x}\|^2 + \eta^2}{s_m}.$$

Finally the first term is bounded by

$$\ln(\rho) - \frac{\#S}{2} \ln[2\pi s_M] - \frac{\|\mathbf{x}\|^2 + \eta^2}{s_m} \leq \ln \left[\sum_{k=1}^K p_k \Phi(\mathbf{x}^S | \mu_k, \Sigma_{k(m)}) \mathbb{I}_{z=k} \right] \leq -\frac{\#S}{2} \ln[2\pi s_m]. \quad (20)$$

Study of the second term:

The second term is expressed as follows:

$$\begin{aligned}
\ln [\Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)})] &= \ln \left[|2\pi \Omega_{(r)}|^{-1/2} \exp \left(-\frac{1}{2} \|\mathbf{x}^U - a - \mathbf{x}^R \beta\|_{\Omega_{(r)}^{-1}}^2 \right) \right] \\
&= -\frac{\#U}{2} \ln[2\pi] - \frac{1}{2} \ln[|\Omega_{(r)}|] - \frac{1}{2} \|\mathbf{x}^U - a - \mathbf{x}^R \beta\|_{\Omega_{(r)}^{-1}}^2.
\end{aligned}$$

Using Lemma 5, the following upper bound is found

$$\ln [\Phi(\mathbf{x}^U | a + \mathbf{x}^R \beta, \Omega_{(r)})] \leq -\frac{\#U}{2} \ln[2\pi s_m].$$

According to Lemma 5, $|\Omega_{(r)}| \leq s_M^{\#U}$ and $\|\mathbf{x}^U - a - \mathbf{x}^R \beta\|_{\Omega_{(r)}^{-1}}^2 \leq s_m^{-1} \|\mathbf{x}^U - a - \mathbf{x}^R \beta\|^2$. In addition,

$$\begin{aligned}
\|\mathbf{x}^U - a - \mathbf{x}^R \beta\|^2 &\leq 2(\|\mathbf{x}^U\|^2 + \|a + \mathbf{x}^R \beta\|^2) \\
&\leq 2(\|\mathbf{x}^U\|^2 + \|a\|^2 + \|\beta\|^2 \|\mathbf{x}^R\|^2) \\
&\leq 2(\|\mathbf{x}^U\|^2 + \eta^2 [1 + \|\mathbf{x}^R\|^2])
\end{aligned}$$

because $a \in \mathcal{B}(\eta, 1, \#U)$ and $\beta \in \mathcal{B}(\eta, \#R, \#U)$. Moreover, $\|\mathbf{x}^U\|^2 \leq \|\mathbf{x}\|^2$ and $\|\mathbf{x}^R\|^2 \leq \|\mathbf{x}\|^2$ hence

$$\|\mathbf{x}^U - a - \mathbf{x}^R\beta\|^2 \leq 2([1 + \eta^2]\|\mathbf{x}\|^2 + \eta^2).$$

Then a lower bound of $\ln[\Phi(\mathbf{x}^U | a + \mathbf{x}^R\beta, \Omega_{(r)})]$ is

$$\ln [\Phi(\mathbf{x}^U | a + \mathbf{x}^R\beta, \Omega_{(r)})] \geq -\frac{\#U}{2} \ln[2\pi s_M] - \frac{\eta^2}{s_m} - \frac{1 + \eta^2}{s_m} \|\mathbf{x}\|^2.$$

Finally the second term is bounded by

$$-\frac{\#U}{2} \ln[2\pi s_M] - \frac{\eta^2}{s_m} - \frac{1 + \eta^2}{s_m} \|\mathbf{x}\|^2 \leq \ln [\Phi(\mathbf{x}^U | a + \mathbf{x}^R\beta, \Omega_{(r)})] \leq -\frac{\#U}{2} \ln[2\pi s_m]. \quad (21)$$

The third term

$$\begin{aligned} \ln [\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})] &= \ln \left[|2\pi\tau_{(l)}|^{-1/2} \exp \left(-\frac{1}{2} \|\mathbf{x}^W - \gamma\|_{\tau_{(l)}^{-1}}^2 \right) \right] \\ &= -\frac{\#W}{2} \ln[2\pi] - \frac{1}{2} \ln[|\tau_{(l)}|] - \frac{1}{2} \|\mathbf{x}^W - \gamma\|_{\tau_{(l)}^{-1}}^2, \end{aligned}$$

can be upper bounded by

$$\ln [\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})] \leq -\frac{\#W}{2} \ln[2\pi s_m],$$

from Lemma 5. According to Lemma 5, $|\tau_{(l)}| \leq s_M^{\#W}$ and

$$\begin{aligned} \|\mathbf{x}^W - \gamma\|_{\tau_{(l)}^{-1}}^2 &\leq s_m^{-1} \|\mathbf{x}^W - \gamma\|^2 \\ &\leq \frac{2}{s_m} (\|\mathbf{x}^W\|^2 + \|\gamma\|^2) \\ &\leq \frac{2}{s_m} (\|\mathbf{x}\|^2 + \eta^2) \end{aligned}$$

because $\gamma \in \mathcal{B}(\eta, \#W)$. Then a lower bound of $\ln[\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})]$ is

$$\ln[\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})] \geq -\frac{\#W}{2} \ln[2\pi s_M] - \frac{(\|\mathbf{x}\|^2 + \eta^2)}{s_m}.$$

Finally the third term is bounded by

$$-\frac{\#W}{2} \ln[2\pi s_M] - \frac{(\|\mathbf{x}\|^2 + \eta^2)}{s_m} \leq \ln [\Phi(\mathbf{x}^W | \gamma, \tau_{(l)})] \leq -\frac{\#W}{2} \ln[2\pi s_m]. \quad (22)$$

Using (20), (21), (22) and $\#S + \#U + \#W = Q$, each function of the family $\mathcal{F}_{(m,r,l,\mathbf{V})}$ is bounded by

$$\ln(\rho) - \frac{Q}{2} \ln[2\pi s_M] - \frac{3(\|\mathbf{x}\|^2 + \eta^2)}{s_m} - \frac{\eta^2 \|\mathbf{x}\|^2}{s_m} \leq \ln[f(\mathbf{x}, z | m, r, l, \mathbf{V}, \theta)] \leq -\frac{Q}{2} \ln[2\pi s_m].$$

Thus, for all $\theta \in \Theta'_{(m,r,l,\mathbf{V})}$ and all $\mathbf{x} \in \mathbb{R}^Q$, $z \in \{1, \dots, K\}$, $|\ln[f(\mathbf{x}, z | m, r, l, \mathbf{V}, \theta)]| \leq C_1(s_m, s_M, Q, \eta, \rho) + C_2(\eta, s_m) \|\mathbf{x}\|^2$ defining the envelop function F , where $C_1(s_m, s_M, Q, \eta, \rho)$

and $C_2(\eta, s_m)$ are two positive constants. To verify that F is h -integrable, we have to show that $\int \|\mathbf{x}\|^2 h(\mathbf{x}, z) d(\mathbf{x}, z) < \infty$:

$$\begin{aligned} \int \|\mathbf{x}\|^2 h(\mathbf{x}, z) d(\mathbf{x}, z) &= \int \|\mathbf{x}\|^2 f(\mathbf{x}, z | m_0, r_0, l_0, \mathbf{V}_0, \theta^*) d(\mathbf{x}, z) \\ &= \sum_{k=1}^K \left[\int \|\mathbf{x}\|^2 f(\mathbf{x} | z = k, m_0, r_0, l_0, \mathbf{V}_0, \theta^*) d\mathbf{x} \right] f(z = k | m_0, r_0, l_0, \mathbf{V}_0, \theta^*) \end{aligned}$$

The integral

$$(1) := \int \|\mathbf{x}\|^2 f(\mathbf{x} | z = k, m_0, r_0, l_0, \mathbf{V}_0, \theta^*) d\mathbf{x}$$

is bounded by

$$\begin{aligned} (1) &= \int \|\mathbf{x}\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \Phi(\mathbf{x}^{U_0} | a^* + \mathbf{x}^{R_0} \beta^*, \Omega_{(r_0)}^*) \Phi(\mathbf{x}^{W_0} | \gamma^*, \tau_{(l_0)}^*) d\mathbf{x}^{W_0} d\mathbf{x}^{U_0} d\mathbf{x}^{S_0} \\ &= \int \|\mathbf{x}^{S_0}\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &+ \int \|\mathbf{x}^{U_0}\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \Phi(\mathbf{x}^{U_0} | a^* + \mathbf{x}^{R_0} \beta^*, \Omega_{(r_0)}^*) d\mathbf{x}^{U_0} d\mathbf{x}^{S_0} \\ &+ \int \|\mathbf{x}^{W_0}\|^2 \Phi(\mathbf{x}^{W_0} | \gamma^*, \tau_{(l_0)}^*) d\mathbf{x}^{W_0} \\ &\leq \int \|\mathbf{x}^{S_0}\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &+ \int 2\|a^* + \mathbf{x}^{R_0} \beta^*\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &+ \int 2\|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \Phi(\mathbf{x}^{U_0} | a^* + \mathbf{x}^{R_0} \beta^*, \Omega_{(r_0)}^*) d\mathbf{x}^{U_0} d\mathbf{x}^{S_0} \\ &+ \int \|\mathbf{x}^{W_0}\|^2 \Phi(\mathbf{x}^{W_0} | \gamma^*, \tau_{(l_0)}^*) d\mathbf{x}^{W_0} \\ &\leq A_1 + A_2 + A_3 + A_4 \end{aligned} \tag{23}$$

The first integral

$$\begin{aligned} A_1 &= \int \|\mathbf{x}^{S_0}\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &\leq [2\|\mu_k\|^2 + 2 \operatorname{tr}(\Sigma_{k(m_0)})] \end{aligned}$$

according to Lemma 6. Thus, from Lemma 5,

$$A_1 \leq 2\eta^2 + 2s_M \# S_0.$$

The second integral is upper bounded by

$$\begin{aligned} A_2 &= \int \|a^* + \mathbf{x}^{R_0} \beta^*\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &\leq \int \eta^2 (1 + \|\mathbf{x}^{S_0}\|^2) \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \\ &\leq \eta^2 \int \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} + \eta^2 A_1 \\ &\leq \eta^2 + \eta^2 [2\eta^2 + 2s_M \# S_0]. \end{aligned}$$

The third integral can be written

$$\begin{aligned}
A_3 &= \int \|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2 \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \Phi(\mathbf{x}^{U_0} | a^* + \mathbf{x}^{R_0} \beta^*, \Omega_{(r_0)}^*) d\mathbf{x}^{U_0} d\mathbf{x}^{S_0} \\
&= \int \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \int \|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2 |2\pi \Omega_{(r_0)}^*|^{-\frac{1}{2}} \\
&\quad \exp \left[-\frac{\|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|_{\Omega_{(r_0)}^*}^2}{2} \right] d\mathbf{x}^{U_0} d\mathbf{x}^{S_0} \\
&\leq \int \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) \int \|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2 (2\pi s_M)^{-\frac{\#U_0}{2}} \\
&\quad \exp \left[-\frac{\|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2}{2s_M} \right] d\mathbf{x}^{U_0} d\mathbf{x}^{S_0}
\end{aligned}$$

because $|\Omega_{(r_0)}^*|^{-1/2} \leq s_M^{-\#U_0/2}$ and $\|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|_{\Omega_{(r_0)}^*}^2 \geq s_M^{-1} \|\mathbf{x}^{U_0} - a^* - \mathbf{x}^{R_0} \beta^*\|^2$ according to Lemma 5. Thus, from Lemma 6,

$$\begin{aligned}
A_3 &\leq \int \Phi(\mathbf{x}^{S_0} | \mu_k^*, \Sigma_{k(m_0)}^*) d\mathbf{x}^{S_0} \times \int \|u\|^2 \Phi(u|0, bI_{\#U_0}) du \\
&= s_M \#U_0.
\end{aligned}$$

The fourth term

$$\begin{aligned}
\int \|\mathbf{x}^{W_0}\|^2 \Phi(\mathbf{x}^{W_0} | \gamma^*, \tau_{(l_0)}^*) d\mathbf{x}^{W_0} &= \int \|\mathbf{x}^{W_0}\|^2 \Phi(\mathbf{x}^{W_0} | \gamma^*, \tau_{(l_0)}^*) d\mathbf{x}^{W_0} \\
&\leq 2[\|\gamma^*\|^2 + \text{tr}(\tau^*)] \\
&\leq 2(\eta^2 + \#W_0 s_M)
\end{aligned}$$

according to Lemma 6.

Thus turning back to Inequality (23), the integral $\int \|\mathbf{x}\|^2 h(\mathbf{x}, z) d\mathbf{x} dz < \infty$ and finally F is h -integrable. Since $\ln(h) \in \mathcal{F}_{(m_0, r_0, l_0, \mathbf{v}_0)}$, it implies that $\mathbb{E}[|\ln h(X, Z)|] < \infty$ and the law of large numbers can be applied to end the proof. \square

Proposition 4.

Assume that

1. (Y_1, \dots, Y_n) is a n -sample with unknown density h .
2. Θ is a compact metric space.
3. $\theta \in \Theta \mapsto \ln[f(\mathbf{y}|\theta)]$ is continuous for every $\mathbf{y} \in \mathbb{R}^Q$.
4. F is an envelope function of $\mathcal{F} := \{\ln[f(\cdot|\theta)]; \theta \in \Theta\}$ which is h -integrable.
5. $\theta^* = \arg\max_{\theta \in \Theta} KL[h, f(\cdot|\theta)]$
6. $\hat{\theta} = \arg\max_{\theta \in \Theta} \sum_{i=1}^n f(Y_i|\theta)$.

Then $\frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\hat{\theta})] \xrightarrow[n \rightarrow \infty]{P} \mathbb{E}_Y[\ln f(Y|\theta^*)]$.

Proof. We consider the following inequality

$$\begin{aligned} \left| \mathbb{E}_Y[\ln f(Y|\theta^*)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\hat{\theta})] \right| &\leq \left| \mathbb{E}_Y[\ln f(Y|\theta^*)] - \mathbb{E}_Y[\ln f(Y|\hat{\theta})] \right| \\ &\quad + \sup_{\theta \in \Theta} \left| \mathbb{E}_Y[\ln f(Y|\theta)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\theta)] \right|. \end{aligned}$$

According to the definition of θ^* , $\mathbb{E}_Y[\ln(f(Y|\theta^*))] - \mathbb{E}_Y[\ln(f(Y|\hat{\theta}_n))] \geq 0$, thus

$$\begin{aligned} \left| \mathbb{E}_Y[\ln f(Y|\theta^*)] - \mathbb{E}_Y[\ln f(Y|\hat{\theta})] \right| &= \mathbb{E}_Y[\ln f(Y|\theta^*)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\theta^*)] \\ &\quad + \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\theta^*)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\hat{\theta})] \\ &\quad + \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\hat{\theta})] - \mathbb{E}_Y[\ln f(Y|\hat{\theta})] \\ &\leq 2 \sup_{\theta \in \Theta} \left| \mathbb{E}_Y[\ln f(Y|\theta)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\theta)] \right|. \end{aligned}$$

According to Example 19.8 in van der Vaart (1998), the bracketing numbers of \mathcal{F} are finite under the assumptions. Hence, using Theorem 19.4 in van der Vaart (1998), \mathcal{F} is P-Glivenko-Cantelli. Thus $\sup_{\theta \in \Theta} \left| \mathbb{E}_Y[\ln f(Y|\theta)] - \frac{1}{n} \sum_{i=1}^n \ln[f(Y_i|\theta)] \right| \xrightarrow[n \rightarrow \infty]{P} 0$, which concludes the proof. \square

Lemma 5. Let $\Sigma \in \mathcal{D}_r$ where \mathcal{D}_r is defined in (H2). Then

1. $s_m^r \leq |\Sigma| \leq s_M^r$ and $\text{tr}(\Sigma) \leq s_M r$
2. $\forall \mathbf{x} \in \mathbb{R}^r, s_M^{-1} \|\mathbf{x}\|^2 \leq \|\mathbf{x}\|_{\Sigma^{-1}}^2 \leq s_m^{-1} \|\mathbf{x}\|^2$

Proof. The proof is based on the eigenvalue decomposition of the variance matrix Σ and the bounded constraint on the eigenvalues because $\Sigma \in \mathcal{D}_r$. \square

Lemma 6.

Let $\Phi(\cdot|\mu, \Sigma)$ be the density of the multivariate Gaussian distribution $\mathcal{N}_r(\mu, \Sigma)$. Then

1. $\int \|\mathbf{x}\|^2 \Phi(\mathbf{x}|0, \Sigma) d\mathbf{x} = \text{tr}(\Sigma)$
2. $\int \|\mathbf{x}\|^2 \Phi(\mathbf{x}|\mu, \Sigma) d\mathbf{x} \leq 2 [\|\mu\|^2 + \text{tr}(\Sigma)]$

Proof. The first result is a classical property of multivariate Gaussian densities. The second result is deduced from the first one using the triangle inequality. \square

Lemma 7.

Let A and B be two real random variables,

$$\forall \epsilon \in \mathbb{R}, P(A + B \leq 0) \leq P(A \leq \epsilon) + P(-B > \epsilon).$$

D The forward variable selection in regression

The following algorithm allows us to determine the subset $R[u]$ of variables among S required to explain \mathbf{x}^u with a linear regression, u being a set of redundant variables. The model comparison is performed with criterion BIC_{reg} defined in (4). The algorithm is making use of the inclusion and exclusion steps now described.

Initialisation $R[u] = \emptyset$, $j_E = \emptyset$ and $j_I = \emptyset$.

Inclusion step For all j in $S \setminus R[u]$, compute

$$\text{B}_{\text{diffreg}}(j) = \text{BIC}_{\text{reg}}(\mathbf{x}^u | r, \mathbf{x}^{R[u] \cup j}) - \text{BIC}_{\text{reg}}(\mathbf{x}^u | r, \mathbf{x}^{R[u]}).$$

Then, compute $j_I = \underset{j \in S \setminus R[u]}{\text{argmax}} \text{B}_{\text{diffreg}}(j)$.

- If $\text{B}_{\text{diffreg}}(j_I) > 0$,
 - if $j_I = j_E$, stop
 - otherwise, $R[u] = R[u] \cup j_I$ and go to the exclusion step.
- Otherwise, $j_I = \emptyset$. If $j_E \neq \emptyset$, go to the exclusion step and stop otherwise.

Exclusion step For all j in $R[u]$, compute

$$\text{B}_{\text{diffreg}}(j) = \text{BIC}_{\text{reg}}(\mathbf{y}^u | r, \mathbf{y}^{R[u]}) - \text{BIC}_{\text{reg}}(\mathbf{x}^u | r, \mathbf{x}^{R[u] - j}).$$

Then, compute $j_E = \underset{j \in R[u]}{\text{argmin}} \text{B}_{\text{diffreg}}(j)$.

- If $\text{B}_{\text{diffreg}}(j_E) \leq 0$, set $R[u] = R[u] - j_E$ and go to the inclusion step if $j_E \neq j_I$ or stop otherwise.
- otherwise, $j_E = \emptyset$ and go to the inclusion step.

Starting from the inclusion step, the forward variable selection algorithm consists of alternating the inclusion and exclusion steps.

References

- Banfield, J. D. and Raftery, A. E. (1993). Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49(3):803–821.
- Bensmail, H. and Celeux, G. (1996). Regularized Gaussian Discriminant Analysis Through Eignenvalue Decomposition. *Journal of the American Statistical Association*, 91(436):1743–1748.
- Biernacki, C., Celeux, G., Govaert, G., and Langrognnet, F. (2006). Model-based cluster and discriminant analysis with the MIXMOD software. *Computational Statistics and Data Analysis*, 51(2):587–600.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer, New York.

- Celeux, G. and Govaert, G. (1995). Gaussian parsimonious clustering models. *Pattern Recognition*, 28(5):781–793.
- Fraley, C. and Raftery, A. E. (2003). Enhanced software for model-based clustering, density estimation, and discriminant analysis: MCLUST. *Journal of Classification*, 20(2):263–286.
- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., and Lander, E. S. (1999). Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring. 286:531–537.
- Guyon, I. and Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine Learning research*, 3:1157–1182.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. Springer, New York, Second Edition.
- Krishnapuram, B., Carin, L., and Hartemink, A. (2004). *Gene expression analysis: joint feature selection and classifier design*, *Kernel Methods in Computational Biology*. MIT Press, Cambridge, MA.
- Mary-Huard, T. and Robin, S. (2009). Tailored aggregation for classification. *IEEE transactions on pattern analysis and machine intelligence*, 31:2098–2105.
- Mary-Huard, T., Robin, S., and Daudin, J.-J. (2007). A penalized criterion for variable selection in classification. *Journal of Multivariate Analysis*, 98:695–705.
- Maugis, C., Celeux, G., and Martin-Magniette, M.-L. (2009a). Variable Selection for Clustering with Gaussian Mixture Models. *Biometrics*, 65:701–709.
- Maugis, C., Celeux, G., and Martin-Magniette, M.-L. (2009b). Variable selection in model-based clustering: A general variable role modeling. *Computational Statistics and Data Analysis*, 53:3872–3882.
- McLachlan, G. (1992). *Discriminant Analysis and Statistical Pattern Analysis*. Wiley-Interscience, New York.
- Murphy, B. T., Raftery, A. E., and Dean, N. (2010). Variable Selection and Updating in Model-Based Discriminant Analysis for High-Dimensional Data with Food Authenticity Applications. *Annals of Applied Statistics*, 4. To appear.
- Raftery, A. E. and Dean, N. (2006). Variable Selection for Model-Based Clustering. *Journal of the American Statistical Association*, 101(473):168–178.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464.
- Su, Y., Murali, T., Pavlovic, V., Schaffer, M., and Kasif, S. (2003). RankGene: identification of diagnostic genes based on expression data. *Bioinformatics*, 19:1578–1579.

- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press, Cambridge.
- Yang, A. J. and Xin-Yuan, S. (2010). Bayesian variable selection for disease classification using gene expression data. *Bioinformatics*, 26(2):215–222.
- Young, D. M. and Odell, P. L. (1986). Feature-subset selection for statistical classification problems involving unequal covariance matrices. *Communication in Statistics-Theory and Methods*, 15:137–157.
- Zhang, Q. and Wang, H. (2008). A BIC Criterion for Gaussian Mixture Model Selection with Application in Discriminant Analysis. Technical report, Guanghua School of Management, Peking University.



Centre de recherche INRIA Saclay – Île-de-France
Parc Orsay Université - ZAC des Vignes
4, rue Jacques Monod - 91893 Orsay Cedex (France)

Centre de recherche INRIA Bordeaux – Sud Ouest : Domaine Universitaire - 351, cours de la Libération - 33405 Talence Cedex
Centre de recherche INRIA Grenoble – Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier
Centre de recherche INRIA Lille – Nord Europe : Parc Scientifique de la Haute Borne - 40, avenue Halley - 59650 Villeneuve d'Ascq
Centre de recherche INRIA Nancy – Grand Est : LORIA, Technopôle de Nancy-Brabois - Campus scientifique
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex
Centre de recherche INRIA Paris – Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex
Centre de recherche INRIA Rennes – Bretagne Atlantique : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex
Centre de recherche INRIA Sophia Antipolis – Méditerranée : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex

Éditeur
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399